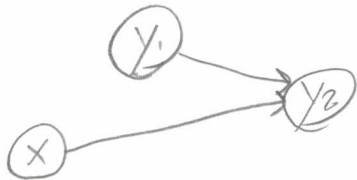


Non couple = <sup>Corona</sup> $X = \text{disease 1 (binary)}$  Heart Attacks $Y_1 = \text{disease 2 (binary)}$  $Y_2 = \text{hospitalization (binary)}$ 

two downstream responses

Here's the DAG:

(E)



We are interested in testing if disease 1 is related to disease 2.

There's no causal arrow from  $X$  to  $Y_1 \Rightarrow$  no relationship.If we run the logistic regression  $\ln\left(\frac{P(Y=1)}{P(Y=0)}\right) = b_0 + b_1 X \Rightarrow b_1 \approx 0$  and insignificant.However, if we run the logistic regression only on  $\mathbb{D}[Y_2=1]$  $\ln\left(\frac{P(Y=1)}{P(Y=0)}\right) = b_0 + b_1 X \Rightarrow b_1 > 0$  and significant Berkson's Paradox

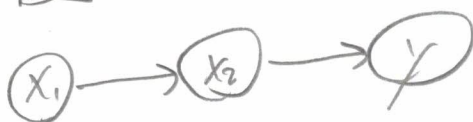
What happens? Since both diseases are necessary to put someone in the hospital, if you're in the hospital, the diseases are suddenly related "collider bias" (1976)

$\Rightarrow$  This illustrates the pitfall of not using simple random samples. In order to estimate both a causal effect, you must sample from the whole population!

$\Rightarrow$  Is the correlation between  $X$  and  $Y_1$  valid for subjects who  $Y_2=1$ ?

Yes! But it is not causal. It is still useful for prediction!

(F)



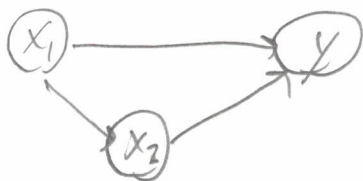
~~X~~: lung cancer  
 $X_1$ : smoking  
 $X_2$ : DNA damage

Is  $X_1$  causal? Yes.  $\hat{y} = b_0 + b_1 x_1$  will  $b_1$  be unbiased? Yes!

What about  $\hat{y} = b_0 + b_1 x_1 + b_2 x_2$ . Will  $b_1$  be unbiased? No!  
 The path is blocked!

(G)

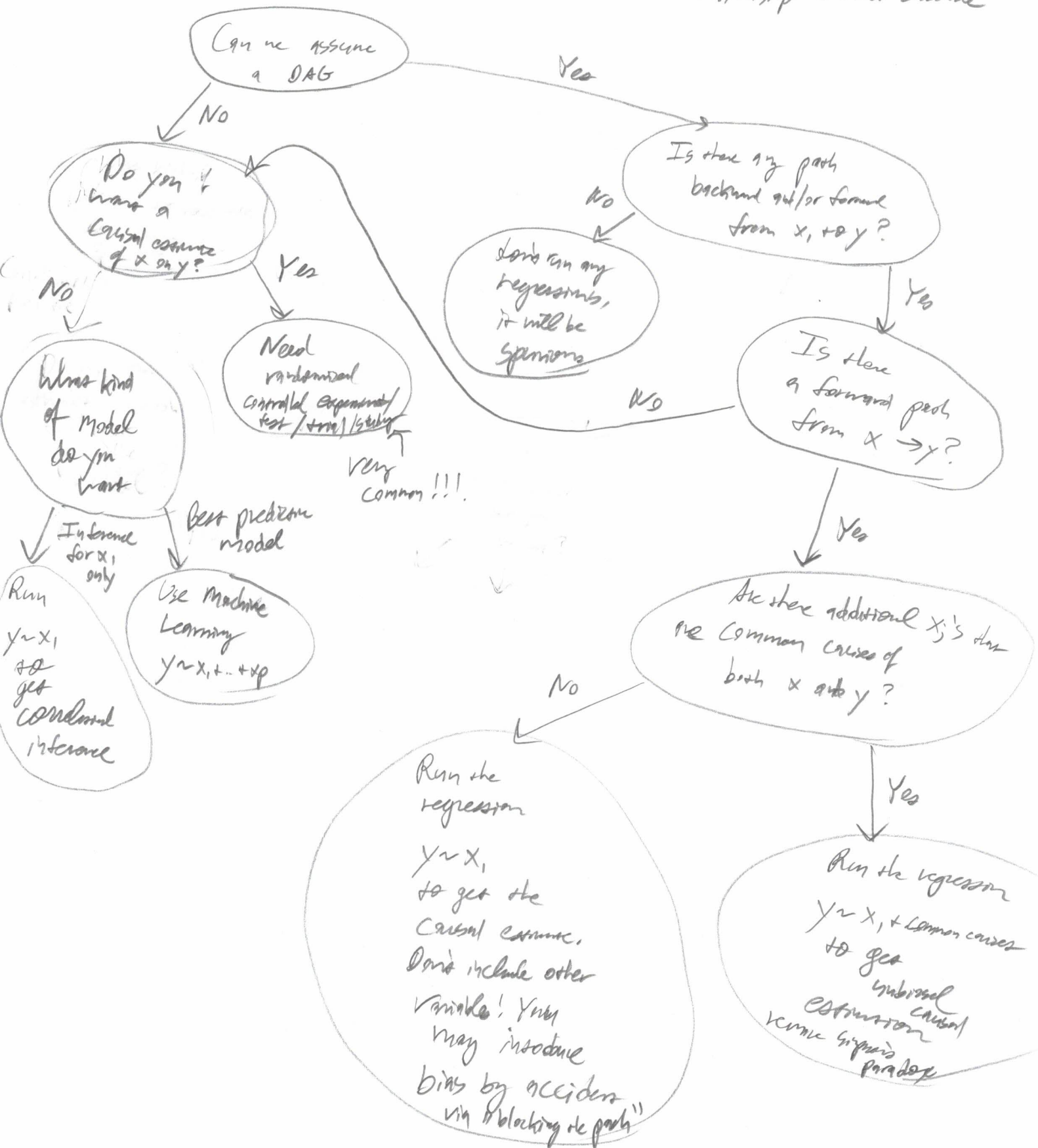
This is now the opposite of Simpson's Paradox



Is  $X_1$  causal? Yes.  $\hat{y} = b_0 + b_1 x_1$  will be unbiased? Yes!

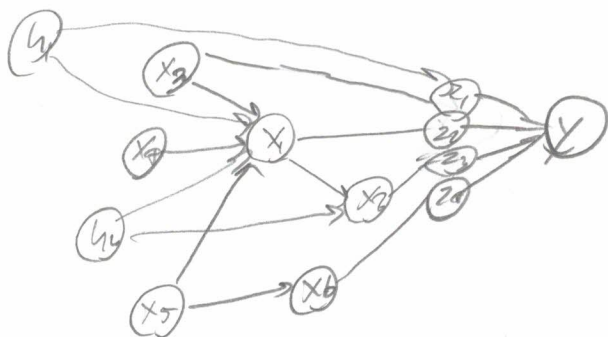
What about  $\hat{y} = b_0 + b_1 x_1 + b_2 x_2$ .  $b_1$  is now biased. The path is partially blocked!

Let's sum up our scenarios:  $Y$  = response,  $X_1$  = variable of interest  
 $X_2, \dots, X_p$  = other variable



What is a controlled experiment / test / trial / study?

Consider the real but unknown DAG:

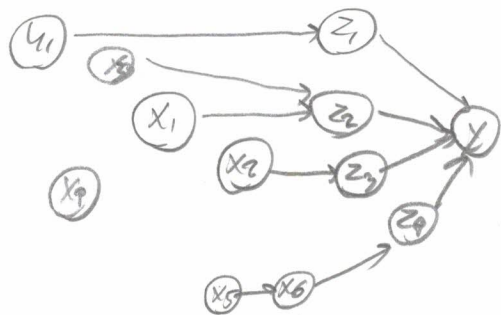


If we run the regression of  $y$  on all observed covariates  $\Rightarrow$  wrong

( . . . . . ) on  $x_1 \Rightarrow$  wrong

( . . . . . ) on any subset of observed covariates  $\Rightarrow$  wrong

If we now do "surgery" on reality and delete all ~~ignores~~ <sup>ignores</sup>  $x_1$ , we have



This deletion of ~~ignores~~ <sup>ignores</sup> now clears the path for  $\hat{y} = b_0 + b_1 x$  being unbiased for the causal effect of  $x_1$  on  $y$ .

How to delete all ~~ignores~~ <sup>ignores</sup>? Control / manipulate values of  $x_1$  yourself.

If it's not manipulated it is called an "observational study" as the data is merely sampled and observed "in the wild".

Is this always possible?  $x_i = \text{Smoking}$ ,  $y = \text{lung cancer}$  <sup>gets</sup>

Can you force people to smoke? Can you force people not to smoke? No

Is this always ethical?  $x_i = \text{experimental and dangerous drug}$  <sup>taking</sup>  
 $y = \text{cancer}$  <sup>gets</sup>

No!

If possible, is it practical and affordable?

$x_i = \text{perfect lab-grown hygienic diet at } \$1000/\text{meal}$  <sup>has</sup>,  $y = \text{survival}$

But, if possible, affordable and ethical, controlled trials can provide causal estimates. It is not a panacea; all problems from 3&1 are still present: type I / type II errors / power.

Now... the details of experiments. Consider  $x_i$  to be a binary variable and call it  $w$  so that  $\vec{w} \in \{0,1\}^n$

$$D = \left( \underbrace{\begin{bmatrix} \vec{w} & \vec{x}_1 & \dots & \vec{x}_p \\ \downarrow & \downarrow & & \downarrow \end{bmatrix}}_X, \vec{y} \right)$$

$\vec{w}$  is said to be a "binary manipulation", "binary treatment".

The subjects  $\{i: w_i=0\}$  and  $\{i: w_i=1\}$  are called the "two arms" of the study. Typically, the  $w=0$  is a "placebo" or "no treatment" or "business as usual" and the  $w=1$  is the "new treatment" in which case if  $\vec{y} = b_0 + b_1 w$ ,  $b_1$  is the causal estimate of the new treatment over no treatment on  $y$ . Call "pill-placebo" controlled trial

18

denotes

If  $w=0$  denotes treatment A and  $w=1$  denotes treatment B, then  $b_1$  is the causal estimate of the difference of B over A. This is sometimes called a "comparative study" or an "A-B test".

---

The vector  $\vec{w}$  is called the "assignment" or "allocation" since you are assigning/forcing/manipulating subject  $i$  to receive  $w_i$ . How do we generate  $\vec{w}$ 's? There are  $2^n$  possibilities! The means of generating  $\vec{w}$ 's is called the "design". What design is "best"? Difficult problem! My whole research program of the past 10yr focuses on it! To measure "best" we need to focus on specific settings and make assumptions.